

BLOCK KACZMARZ METHOD WITH INEQUALITIES

J. BRISKMAN AND D. NEEDELL

ABSTRACT. The randomized Kaczmarz method is an iterative algorithm that solves systems of linear equations. Recently, the randomized method was extended to systems of equalities and inequalities by Leventhal and Lewis. Even more recently, Needell and Tropp provided an analysis of a block version of this randomized method for systems of linear equations. This paper considers the use of a block type method for systems of mixed equalities and inequalities, bridging these two bodies of work. We show that utilizing a matrix paving over the equalities of the system can lead to significantly improved convergence, and prove a linear convergence rate as in the standard block method. We also demonstrate that using blocks of inequalities offers similar improvement only when the system satisfies a certain geometric property. We support the theoretical analysis with several experimental results.

1. INTRODUCTION

The Kaczmarz method [18] is an iterative algorithm for solving linear systems of equations. It is usually applied to large-scale overdetermined systems because of its simplicity and speed (but also converges in the underdetermined case to the least-norm solution under appropriate initial conditions). Each iteration projects onto the solution space corresponding to one row in the system, in a sequential fashion. Strohmer and Vershynin prove that when the rows are selected from a certain random distribution rather than sequentially, that the randomized method converges to the solution at a linear rate [31]. The method has been applied to fields including image reconstruction, digital signal processing, and computer tomography [30, 10, 21, 11]. Leventhal and Lewis modify the randomized Kaczmarz method to apply to systems of linear equalities and inequalities [19], thereby extending results on the standard method in this setting (see e.g. [5] and references therein). Unlike the traditional randomized algorithm which enforces a single constraint at each iteration, the block Kaczmarz approach recently analyzed by Needell and Tropp [24] enforces multiple constraints simultaneously and thus offers computational advantages. Here we demonstrate convergence for a system of linear equalities and inequalities by combining a randomized block Kaczmarz method for the equalities with a randomized Kaczmarz algorithm for the inequalities. These results indicate that the block Kaczmarz method can be used for a system of equalities and inequalities, and in some cases may quicken convergence. We also consider the case of utilizing blocking in both the equalities and inequalities, although this can be detrimental unless the geometry of the system meets certain conditions.

1.1. Model and Notation. We consider a linear system

$$\mathbf{A}\mathbf{x} = \mathbf{b}, \tag{1.1}$$

where \mathbf{A} is a real (or complex) $n \times d$ matrix, typically with $n \gg d$.

The ℓ_p vector norm for $p \in [1, \infty]$ is denoted $\|\cdot\|_p$, while $\|\cdot\|$ is the spectral norm and $\|\cdot\|_F$ refers to the Frobenius norm. For an $n \times d$ matrix \mathbf{A} , the singular values are arranged in decreasing order and we write

$$\sigma_{\max}(\mathbf{A}) \stackrel{\text{def}}{=} \sigma_1(\mathbf{A}) \geq \sigma_2(\mathbf{A}) \geq \cdots \geq \sigma_d(\mathbf{A}) \stackrel{\text{def}}{=} \sigma_{\min}(\mathbf{A}).$$

We define the eigenvalues $\lambda_{\min}(\mathbf{A}), \dots, \lambda_{\max}(\mathbf{A})$ of a matrix analogously. For convenience we will assume that each row \mathbf{a}_i of \mathbf{A} has unit norm, $\|\mathbf{a}_i\|_2 = 1$, and we call such matrices *standardized*.

We define the usual condition number

$$\kappa(\mathbf{A}) \stackrel{\text{def}}{=} \sigma_{\max}(\mathbf{A}) / \sigma_{\min}(\mathbf{A}),$$

and write the Moore-Penrose pseudoinverse of matrix \mathbf{A} by \mathbf{A}^\dagger . Recall that for a matrix \mathbf{A} with full row rank, the pseudoinverse is obtained by $\mathbf{A}^\dagger \stackrel{\text{def}}{=} \mathbf{A}^* (\mathbf{A} \mathbf{A}^*)^{-1}$.

Now we consider a system of linear equalities and inequalities and denote by S its non-empty set of feasible solutions. We thus consider the matrix \mathbf{A} whose rows can be arranged such that

$$\mathbf{A} = \begin{bmatrix} \mathbf{A}_= \\ \mathbf{A}_\leq \end{bmatrix}, \quad (1.2)$$

and we will write $I_=$ and I_\leq to denote the row indices of $\mathbf{A}_=$ and \mathbf{A}_\leq , respectively. Therefore, we ask that

$$\langle \mathbf{a}_i, \mathbf{x} \rangle \leq \mathbf{b}_i \quad (i \in I_\leq) \quad \text{and} \quad \langle \mathbf{a}_i, \mathbf{x} \rangle = \mathbf{b}_i \quad (i \in I_=) \quad (1.3)$$

We will assume that the set of rows $\{1, 2, \dots, n\}$ is partitioned such that the first n_e rows correspond to equalities, and the remaining $n_i = n - n_e$ rows to inequalities. Thus $\mathbf{A}_=$ is an $n_e \times d$ matrix and \mathbf{A}_\leq is $n_i \times d$.

The error bound for this system of linear inequalities uses the function $e : \mathbf{R}^n \rightarrow \mathbf{R}^n$ defined as in [19] by

$$e(y)_i = \begin{cases} y_i^+ & \text{for } i \in I_\leq \\ y_i & \text{for } i \in I_= \end{cases}$$

where the positive part is defined as $x^+ \stackrel{\text{def}}{=} \max(x, 0)$.

1.2. Details of Kaczmarz. The simple Kaczmarz method is an iterative algorithm that approximates a least-squares minimizer \mathbf{x}_\star to the problem in (1.1). It takes an arbitrary initial approximation \mathbf{x}_0 , and at each iteration j the current iterate is projected orthogonally onto the solution hyperplane $\{\langle \mathbf{a}_i, \mathbf{x} \rangle = \mathbf{b}_i\}$, using the update rule

$$\mathbf{x}_{j+1} = \mathbf{x}_j + \frac{\mathbf{b}_i - \langle \mathbf{a}_i, \mathbf{x}_j \rangle}{\|\mathbf{a}_i\|_2^2} \mathbf{a}_i \quad (1.4)$$

where $i = j \bmod n + 1$ [18]. With an unfortunate ordering of the rows, this method as-is can produce very slow convergence. However, it has been well known that using randomized selection often eliminates this effect [13, 15]. The randomized Kaczmarz method put forth by Strohmer and Vershynin [31] uses a random selection method for the selection of row i such that each row i is selected with probability proportional to $\|\mathbf{a}_i\|_2^2$. This randomization provides an algorithm that is both simple to analyze and enforce in many cases. In this paper we assume each row has unit norm, so each row is selected uniformly at random from $\{1, \dots, n\}$ in the simple randomized Kaczmarz approach.¹ Strohmer and Vershynin prove a linear rate of convergence for consistent systems that depends on the scaled condition number of \mathbf{A} , and not on the number of equations n [31],

$$\mathbb{E} \|\mathbf{x}_j - \mathbf{x}_\star\|_2^2 \leq \left[1 - \frac{1}{K} \right]^j \|\mathbf{x}_0 - \mathbf{x}_\star\|_2^2, \quad (1.5)$$

where \mathbf{x}_\star is the solution to the consistent system (1.1) and $K = \|\mathbf{A}\|_F^2 / \sigma_{\min}^2(\mathbf{A})$ denotes the scaled condition number. Needell extended this work to the inconsistent case and proves linear convergence to the least-squares solution within some fixed radius [22],

$$\mathbb{E} \|\mathbf{x}_j - \mathbf{x}_\star\|_2^2 \leq \left[1 - \frac{1}{K} \right]^j \|\mathbf{x}_0 - \mathbf{x}_\star\|_2^2 + K \|\mathbf{e}\|_\infty^2,$$

¹This assumption is both for notational convenience, and because the use of matrix pavings discussed below only hold for standardized matrices. In practice, one can employ pre-conditioning on non-standardized systems, or extend the construction of matrix pavings to non-standardized systems [37].

where $\mathbf{e} = \mathbf{A}\mathbf{x}_\star - \mathbf{b}$ denotes the residual vector. Because the Kaczmarz method projects directly onto each solution hyperplane, such a convergence radius is unavoidable without adding a relaxation parameter.

The randomized Kaczmarz method can be adapted to the case of a linear system of equalities and inequalities described in (1.3). Leventhal and Lewis [19] apply the Kaczmarz method to a consistent system of linear equalities and inequalities (here consistent simply means the feasible set S is non-empty). At each iteration j , the previous iterate only projects onto the solution hyperplane if the inequality is not already satisfied. If the inequality is satisfied for row i selected at iteration j ($\mathbf{a}_i^T \mathbf{x} \leq \mathbf{b}_i$), the approximation \mathbf{x}_j is set as \mathbf{x}_{j-1} [19]. The update rule for this algorithm is thus

$$\mathbf{x}_{j+1} = \mathbf{x}_j - \frac{e(\mathbf{a}_i^T \mathbf{x}_j - \mathbf{b}_i)}{\|\mathbf{a}_i\|_2^2} \mathbf{a}_i. \quad (1.6)$$

This algorithm converges linearly in expectation [19], with

$$\mathbb{E}[d(\mathbf{x}_j, S)^2 | \mathbf{x}_{j-1}] \leq d(\mathbf{x}_{j-1}, S)^2 - \frac{\|e(\mathbf{A}\mathbf{x}_{j-1} - \mathbf{b})\|_2^2}{\|\mathbf{A}\|_F^2}.$$

In order to bound the right hand side of this expression, the authors rely on a lemma due to Hoffman [17, 19]. This result states that for any system (1.3) with non-empty solution set S , there exists a constant L independent of \mathbf{b} such that for all \mathbf{x} ,

$$d(\mathbf{x}, S) \leq L \|e(\mathbf{A}\mathbf{x} - \mathbf{b})\|_2. \quad (1.7)$$

When $\mathbf{A}_\perp = \mathbf{A}$ is full column rank, the Hoffman constant is the inverse of the smallest singular value, $L = \sigma_{\min}^{-1}(\mathbf{A})$.

Using this their result becomes

$$\mathbb{E}[d(\mathbf{x}_j, S)^2] \leq \left[1 - \frac{1}{L^2 \|\mathbf{A}\|_F^2}\right]^j \cdot d(\mathbf{x}_0, S)^2, \quad (1.8)$$

which coincides with (1.5) for consistent systems of equalities.

1.3. Block Kaczmarz. A block variant of the randomized Kaczmarz method due to Elfving [9] has been recently analyzed by Needell and Tropp [24] and can improve the convergence rate in certain cases. The block Kaczmarz method first partitions the rows $\{1, \dots, n\}$ into m blocks, denoted τ_1, \dots, τ_m . Instead of selecting one row per iteration as done with the simple Kaczmarz method, the block Kaczmarz algorithm chooses a block uniformly at random at each iteration. Thus the block Kaczmarz method enforces multiple constraints simultaneously. At each iteration, the previous iterate \mathbf{x}_{j-1} is projected onto the solution space to $\mathbf{A}_\tau \mathbf{x} = \mathbf{b}_\tau$, which enforces the set of equations in block τ [24]. \mathbf{A}_τ and \mathbf{b}_τ are written as the row submatrix of \mathbf{A} and the subvector of \mathbf{b} indexed by τ respectively, yielding an iterative rule of

$$\mathbf{x}_j = \mathbf{x}_{j-1} + (\mathbf{A}_\tau)^\dagger (\mathbf{b}_\tau - \mathbf{A}_\tau \mathbf{x}_{j-1}). \quad (1.9)$$

The pseudoinverse used in (1.9) returns the solution to the underdetermined least squares problem for a wide or square row submatrix \mathbf{A}_τ .

Depending on the characteristics of the submatrix \mathbf{A}_τ , the block method can provide better convergence than the simple method. If we assume that the submatrices \mathbf{A}_τ are well conditioned, the additional cost of computing their pseudo-inverse can be overcome by the gain in utilizing block multiplications (see our experiments in Section 4). In fact, if the blocks admit a fast multiply (for example if the matrix is built of DFT or circulant blocks), then the computational cost of the block iteration (1.9) is similar to the cost of the simple update rule in (1.4). Since the convergence depends heavily on the conditioning of each submatrix, one seeks partitions of the rows into blocks for which each block is well-conditioned. The notion of a *row-paving* allows one to do precisely that.

Definition 1.1. We define an (m, β) row paving² of matrix A as a partition $T = \{\tau_1, \dots, \tau_m\}$ of the row indices such that

$$\lambda_{\max}(A_\tau A_\tau^*) \leq \beta \quad \text{for each } \tau \in T.$$

The *size* of the paving, or number of blocks, is m . The value of β is the upper paving bound, which controls the spectral norms of the submatrices. Needell and Tropp [24] show that these parameters determine the performance of the algorithm, with convergence for a consistent system admitting an (m, β) paving given by

$$\mathbb{E} \|x_j - x_\star\|_2^2 \leq \left[1 - \frac{\sigma_{\min}^2(A)}{\beta m} \right]^j \|x_0 - x_\star\|_2^2. \quad (1.10)$$

Therefore the convergence rate depends on the size m and upper bound β ; the algorithm's performance improves with low values of m and β , and large $\sigma_{\min}^2(A)$. The authors also prove convergence for inconsistent systems, with the same convergence rate and convergence radius which depends also on the minimum of all $\lambda_{\min}(A_\tau A_\tau^*)$, see [24] for details.

Surprisingly, every standardized matrix admits a good row paving. The following result is due to [38, 34] which builds off the foundational work of [1, 2].

Proposition 1.2 (Existence of Good Row Pavings). *For any $\delta \in (0, 1)$ and standardized $n \times d$ matrix A , there is a row paving satisfying*

$$m \leq C \cdot \delta^{-2} \|A\|^2 \log(1+n) \quad \text{and} \quad 1 - \delta \leq \beta \leq 1 + \delta.$$

where C is an absolute constant.

Although this is an existential result, there are constructive methods to obtain such pavings, and for certain classes of matrices, they can even be obtained by a random partitioning of the rows [33, 7, 24].

With such a paving in tow, the convergence of (1.10) becomes

$$\mathbb{E} \|x_j - x_\star\|_2^2 \leq \left[1 - \frac{1}{C\kappa^2(A) \log(1+n)} \right]^j \|x_0 - x_\star\|_2^2$$

Although often comparable to the convergence rate for the simple method (1.5), numerical results confirm that the block method offers significant reduction in computation time due to the speed of matrix-vector multiplication (see e.g. [24]).

1.4. Contribution. This paper analyzes the system with matrix described in (1.2) using an algorithm with the block Kaczmarz approach for the equalities given by A_+ and the simple method for the inequalities given by A_- . A paving is created for A_+ , with the inequalities excluded. At each iteration, we select from A_+ with a fixed probability p and from A_- with probability $1 - p$. In the former case, we select a block τ from paving T uniformly at random, and in the latter case we select a row i of A_- uniformly at random. In the case of a block of equalities being selected, the algorithm proceeds by updating x_j using (1.9). When an inequality row is selected, x_j is updated using the rule (1.6). We prove that this method yields linear convergence to the solution set S . We also include a discussion about paving both A_+ and A_- , which identifies a geometric property of the system which allows for improved convergence by utilizing two pavings. We show that when this property is not satisfied, utilizing both pavings can be detrimental to convergence.

1.5. Organization. Section 2 lays out our main result, Theorem 2.1, and provides a proof. We discuss blocking the full matrix in Section 3 and Section 4 explains numerical experiments and results. We conclude with discussion and related work in Section 5.

²The standard definition of a row paving also includes a constant α which serves as a lower bound to the smallest singular value. We ignore that parameter here since it will not be utilized.

2. ANALYSIS OF THE BLOCK KACZMARZ ALGORITHM FOR A SYSTEM OF INEQUALITIES

In this section we analyze the convergence of the described method, which is detailed in Algorithm 2.1.

Algorithm 2.1 Block Kaczmarz Method for a System of Inequalities**Input:**

- Matrix A with dimension $n \times d$
- Right-hand side b with dimension n
- Number of rows representing equalities, n_e , and inequalities, $n_i = n - n_e$
- Partition $T = \{\tau_1, \dots, \tau_m\}$ of the row indices $\{1, \dots, n_e\}$ and paving constant β
- Initial iterate x_0 with dimension d
- Convergence tolerance $\varepsilon > 0$

Output: An estimate \hat{x} to the solution of the system (1.3)

```

j ← 0
repeat
  j ← j + 1
  Draw uniformly at random q from [0, 1]
  if q ≤  $\frac{\beta m}{n_i + \beta m}$ 
    Choose a block  $\tau$  uniformly at random from T
     $x_j \leftarrow x_{j-1} + (A_\tau)^\dagger (b_\tau - A_\tau x_{j-1})$ 
  else
    Choose a row i uniformly at random from  $\{n_e + 1, \dots, n\}$ 
     $x_j \leftarrow x_{j-1} - \frac{e(a_i^T x_{j-1} - b_i)}{\|a_i\|_2^2} a_i$ 
until  $\|e(Ax_j - b)\|_2^2 \leq \varepsilon^2$ 
 $\hat{x} \leftarrow x_j$ 

```

Notice that the probability of selecting a block of A_- is $\frac{\beta m}{n_i + \beta m}$. This quantity corresponds to the relative size of A_- in the system, where the size is measured in terms of the paving quantities βm . This value may be difficult to compute precisely, and the simpler threshold of n_e/n appears to also work well in practice. We provide no evidence that our selection of this threshold is most efficient, nor any more efficient than using one proportional to the number of equality rows n_e . We find that this algorithm yields linear convergence in expectation with a rate that only depends on the number of inequalities n_i , paving size m , and upper bound β .

Our main result is described in Theorem 2.1.

Theorem 2.1 (Convergence). *Let the standardized matrix $A \in \mathbb{R}^{n \times d}$ and $b \in \mathbb{R}^n$ correspond to a system as in (1.2) with the first n_e rows being equalities and the remaining $n_i = n - n_e$ rows being inequalities. Let T be an (m, β) row paving of A_- . Let x_0 be an arbitrary initial estimate and S the non-empty feasible region. Then Algorithm 2.1 satisfies for each iteration $j = 1, 2, 3, \dots$,*

$$\mathbb{E}[d(x_j, S)^2] \leq \left(1 - \frac{1}{L^2(n_i + \beta m)}\right)^j \cdot d(x_0, S)^2,$$

where L is the Hoffman constant (1.7).

Remarks.

1. Note that when there are no block projections, no inequalities, or neither, Theorem 2.1 recovers the

results of the standard randomized Kaczmarz for inequalities [19], the standard randomized block Kaczmarz method [24] or the standard randomized Kaczmarz method [31], respectively. We thus view this result as a completely generalized convergence bound.

2. If we let ρ_s and ρ_b be the convergence rates of the simple and block methods for mixed systems, respectively, then by (1.8) and Theorem 2.1,

$$\rho_s \geq \frac{1}{L^2 n} \quad \text{and} \quad \rho_b \geq \frac{1}{L^2(n_i + \beta m)}.$$

It is evident that our expected convergence rate will be faster per iteration than the simple method when $n_i + \beta m < n$. Since β can be chosen close to 1 and $m < n_e$ is then number of rows in A_- , this holds quite easily.

3. Since a single iteration using a block A_τ in general may cost more than an iteration utilizing a single row, it is more fair to compare per epoch, rather than per iteration. An epoch is typically the minimum number of iterations needed to visit each row of the matrix. When there are inequalities present that are already satisfied in a given iteration, that iteration may make no contribution and cost very little computationally. Thus the notion of epoch may be slightly skewed here, but if we ignore this subtlety the simple method will have approximately n iterations per epoch, compared to $n_i + m$ iterations per epoch with the block method. The approximate per epoch convergence rates can thus be compared as

$$n \cdot \rho_s \geq \frac{1}{L^2} \quad \text{and} \quad (n_i + m) \cdot \rho_b \geq \frac{n_i + m}{L^2(n_i + \beta m)}.$$

This result is similar to that found by Needell and Tropp [24], with the block convergence rate at best equal to that of the simple convergence rate when $\beta = 1$. However, as already noted, the block method is quite advantageous computationally.

Combining the paving result of Prop. 1.2 with Theorem 2.1 yields the following corollary.

Corollary 2.2. *Instate the assumptions and notation of Theorem 2.1 and let A_- be equipped with an (m, β) row-paving as in Proposition 1.2. Then the iterates of Algorithm 2.1 satisfy*

$$\mathbb{E}[d(\mathbf{x}_j, S)^2] \leq \gamma^j \cdot d(\mathbf{x}_0, S)^2,$$

where $\gamma = \left(1 - \frac{1}{L^2(n_i + C\|A_-\|^2 \log(1+n))}\right)$ and C is some absolute constant.

of Theorem 2.1. Fix an iteration j of Algorithm 2.1. We proceed as in [24] and [19]. First, we suppose that $q \leq \frac{\beta m}{n_i + \beta m}$, so that a block τ of equalities is selected this iteration. Then writing P_S as the orthogonal projection onto S , we have $\mathbf{b}_\tau = A_\tau P_S \mathbf{x}_{j-1}$ since $P_S \mathbf{x}_{j-1} \in S$. We then have

$$\begin{aligned} \mathbf{x}_j &= \mathbf{x}_{j-1} + A_\tau^\dagger(\mathbf{b}_\tau - A_\tau \mathbf{x}_{j-1}) \\ &= \mathbf{x}_{j-1} + A_\tau^\dagger(A_\tau P_S \mathbf{x}_{j-1} - A_\tau \mathbf{x}_{j-1}) \\ &= \mathbf{x}_{j-1} + A_\tau^\dagger A_\tau (P_S \mathbf{x}_{j-1} - \mathbf{x}_{j-1}). \end{aligned}$$

Thus,

$$\begin{aligned} \|\mathbf{x}_j - P_S \mathbf{x}_{j-1}\|^2 &= \|\mathbf{x}_{j-1} - P_S \mathbf{x}_{j-1} - A_\tau^\dagger A_\tau (\mathbf{x}_{j-1} - P_S \mathbf{x}_{j-1})\|_2^2 \\ &= \|(\mathbf{I} - A_\tau^\dagger A_\tau)(\mathbf{x}_{j-1} - P_S \mathbf{x}_{j-1})\|_2^2. \end{aligned}$$

Taking expectation (over the choice of the block τ , conditioned on previous choices), and using the fact that $A_\tau^\dagger A_\tau$ is an orthogonal projector, along with the properties of the paving yields

$$\begin{aligned}
& \mathbb{E} \|\mathbf{x}_j - P_S \mathbf{x}_{j-1}\|^2 \\
&= \mathbb{E} \|(I - \mathbf{A}_\tau^\dagger \mathbf{A}_\tau)(\mathbf{x}_{j-1} - P_S \mathbf{x}_{j-1})\|_2^2 \\
&= \|\mathbf{x}_{j-1} - P_S \mathbf{x}_{j-1}\|_2^2 - \mathbb{E} \|\mathbf{A}_\tau^\dagger \mathbf{A}_\tau (\mathbf{x}_{j-1} - P_S \mathbf{x}_{j-1})\|_2^2 \\
&\leq \|\mathbf{x}_{j-1} - P_S \mathbf{x}_{j-1}\|_2^2 - \frac{1}{\beta} \mathbb{E} \|\mathbf{A}_\tau (\mathbf{x}_{j-1} - P_S \mathbf{x}_{j-1})\|_2^2.
\end{aligned}$$

Since $d(\mathbf{x}_{j-1}, S) = \|\mathbf{x}_{j-1} - P_S \mathbf{x}_{j-1}\|_2$ and $d(\mathbf{x}_j, S) \leq \|\mathbf{x}_j - P_S \mathbf{x}_{j-1}\|_2$, this means that

$$\begin{aligned}
\mathbb{E} [d(\mathbf{x}_j, S)^2] &\leq d(\mathbf{x}_{j-1}, S)^2 - \frac{1}{\beta} \mathbb{E} \|\mathbf{A}_\tau (\mathbf{x}_{j-1} - P_S \mathbf{x}_{j-1})\|_2^2 \\
&= d(\mathbf{x}_{j-1}, S)^2 - \frac{1}{\beta m} \sum_{\tau \in T} \|\mathbf{A}_\tau \mathbf{x}_{j-1} - \mathbf{b}_\tau\|_2^2 \\
&= d(\mathbf{x}_{j-1}, S)^2 - \frac{1}{\beta m} \sum_{i \in I_\tau} e(\mathbf{A}_\tau \mathbf{x}_{j-1} - \mathbf{b}_\tau)_i^2. \tag{2.1}
\end{aligned}$$

Next suppose that instead $i \in I_\leq$ is selected. Then since each row \mathbf{a}_i has unit norm,

$$\begin{aligned}
d(\mathbf{x}_j, S)^2 &\leq \|\mathbf{x}_j - P_S \mathbf{x}_{j-1}\|_2^2 \\
&= \|\mathbf{x}_{j-1} - e(\mathbf{A} \mathbf{x}_{j-1} - \mathbf{b})_i \mathbf{a}_i - P_S \mathbf{x}_{j-1}\|_2^2 \\
&= \|\mathbf{x}_{j-1} - P_S \mathbf{x}_{j-1}\|_2^2 + e(\mathbf{A} \mathbf{x}_{j-1} - \mathbf{b})_i^2 \\
&\quad - 2e(\mathbf{A} \mathbf{x}_{j-1} - \mathbf{b})_i \langle \mathbf{a}_i, \mathbf{x}_{j-1} - P_S \mathbf{x}_{j-1} \rangle \\
&\leq d(\mathbf{x}_{j-1}, S)^2 - e(\mathbf{A} \mathbf{x}_{j-1} - \mathbf{b})_i^2,
\end{aligned}$$

where the last line follows from the fact that $\langle \mathbf{a}_i, P_S \mathbf{x}_{j-1} \rangle \leq b_i$ and $e(\mathbf{A} \mathbf{x}_{j-1} - \mathbf{b})_i \geq 0$. Now taking expectation again we have

$$\begin{aligned}
\mathbb{E} [d(\mathbf{x}_j, S)^2] &\leq d(\mathbf{x}_{j-1}, S)^2 - \mathbb{E}(e(\mathbf{A} \mathbf{x}_{j-1} - \mathbf{b})_i^2) \\
&= d(\mathbf{x}_{j-1}, S)^2 - \frac{1}{n_i} \sum_{i \in I_\leq} e(\mathbf{A}_\leq \mathbf{x}_{j-1} - \mathbf{b}_\leq)_i^2.
\end{aligned}$$

Combining these results and letting E_τ and E_\leq denote the events that a block from T and a row from I_\leq is selected, respectively, we have

$$\begin{aligned}
\mathbb{E} [d(\mathbf{x}_j, S)^2] &= p \cdot \mathbb{E}[d(\mathbf{x}_j, S)^2 | E_\tau] + (1-p) \cdot \mathbb{E}[d(\mathbf{x}_j, S)^2 | E_\leq] \\
&\leq p \left[d(\mathbf{x}_{j-1}, S)^2 - \frac{1}{\beta m} \sum_{i \in I_\tau} e(\mathbf{A}_\tau \mathbf{x}_{j-1} - \mathbf{b}_\tau)_i^2 \right] \\
&\quad + (1-p) \left[d(\mathbf{x}_{j-1}, S)^2 - \frac{1}{n_i} \sum_{i \in I_\leq} e(\mathbf{A}_\leq \mathbf{x}_{j-1} - \mathbf{b}_\leq)_i^2 \right] \\
&= d(\mathbf{x}_{j-1}, S)^2 - p \cdot \frac{1}{\beta m} \sum_{i \in I_\tau} e(\mathbf{A}_\tau \mathbf{x}_{j-1} - \mathbf{b}_\tau)_i^2 \\
&\quad - (1-p) \cdot \frac{1}{n_i} \sum_{i \in I_\leq} e(\mathbf{A}_\leq \mathbf{x}_{j-1} - \mathbf{b}_\leq)_i^2.
\end{aligned}$$

Since $p = \frac{\beta m}{n_i + \beta m}$, we have $\frac{1-p}{n_i} = \frac{1}{n_i + \beta m}$ and we can simplify

$$\begin{aligned} \mathbb{E} [d(\mathbf{x}_j, S)^2] &\leq d(\mathbf{x}_{j-1}, S)^2 - \frac{1}{n_i + \beta m} \left[\sum_{i \in I_{=}} e(\mathbf{A}_{=}\mathbf{x}_{j-1} - \mathbf{b}_{=})_i^2 \right. \\ &\quad \left. + \sum_{i \in I_{\leq}} e(\mathbf{A}_{\leq}\mathbf{x}_{j-1} - \mathbf{b}_{\leq})_i^2 \right] \\ &= d(\mathbf{x}_{j-1}, S)^2 - \frac{1}{n_i + \beta m} \|\mathbf{e}(\mathbf{A}\mathbf{x}_{j-1} - \mathbf{b})\|_2^2 \\ &\leq d(\mathbf{x}_{j-1}, S)^2 - \frac{1}{L^2(n_i + \beta m)} \cdot d(\mathbf{x}_{j-1}, S)^2 \\ &= \left[1 - \frac{1}{L^2(n_i + \beta m)} \right] d(\mathbf{x}_{j-1}, S)^2, \end{aligned}$$

where we have utilized the Hoffman bound (1.7) in the second inequality.

Utilizing independence of the random selections and recursing on this relation yields the desired result. \square

3. A DISCUSSION ABOUT BLOCKING INEQUALITIES

It is natural to ask whether one can benefit by blocking both the equalities as above and also the inequalities, as described by Algorithm 3.2. Indeed, Section 4 will show dramatic improvements in computational time when the rows of $\mathbf{A}_{=}$ are paved and block projections as in Algorithm 2.1 are used. So can one benefit even more by paving also the rows of \mathbf{A}_{\leq} ? The answer to this question heavily depends on the structure of the matrix \mathbf{A} .

If we only consider $\mathbf{A}_{=}$, a block projection as in (1.9) enforces all the equations indexed by τ to be satisfied. This is of course desirable when the rows indexed by τ correspond to equalities. Also, if a single inequality corresponding to row i in \mathbf{A}_{\leq} is not satisfied and we perform a single projection as in (1.4), we are again enforcing that inequality to hold with equality. However, this improves the estimation since in this case we know the solution set S lies on the opposite side of the hyperplane $\{\mathbf{x} : \langle \mathbf{x}, \mathbf{a}_i \rangle = b_i\}$ as the current estimation (see Figure 1 (a)). On the other hand, if we employ a block projection as in (1.9) to a

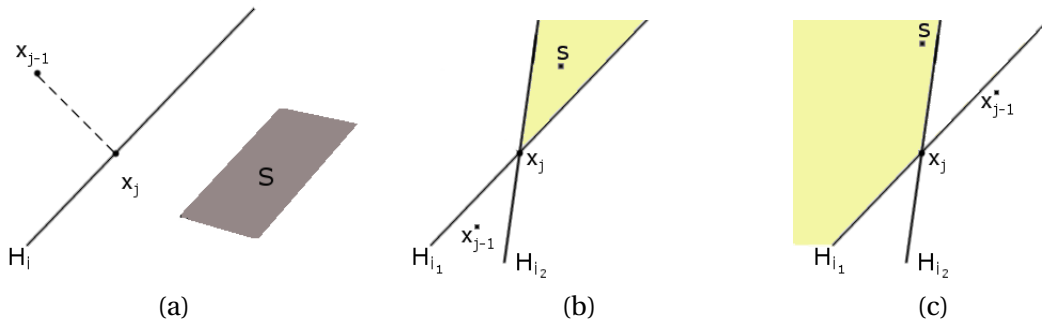


Figure 1 Possible geometries of the system. S denotes solution space (or solution point). Yellow shading denotes regions where inequalities i_1 and i_2 are both satisfied. (a) A single projection onto hyperplane $H_i = \{\mathbf{x} : \langle \mathbf{a}_i, \mathbf{x} \rangle = b_i\}$ provides improved estimation. (b) Block projection onto intersection of hyperplanes also may provide improved estimation. (c) Block projection onto intersection of hyperplanes may provide improved estimation.

set of inequalities indexed by τ which are not satisfied by the current estimation \mathbf{x}_{j-1} then we enforce *all of them to hold with equality simultaneously*. Depending on the geometry of the involved rows, this may

result in an improved estimation or actually one much farther from the solution set. Of course, one might alternatively want to solve the convex program to project onto the intersection of the corresponding half-spaces, but we would like to maintain the efficiency and simplicity of the block Kaczmarz method.

As an illustrative example, Figure 1 (b) and (c) demonstrate two possible scenarios in two dimensions. Here, the solution space is a single point marked S , and we draw two hyperplanes H_{i_1} and H_{i_2} where $H_i = \{\mathbf{x} : \langle \mathbf{a}_i, \mathbf{x} \rangle = b_i\}$. The yellow shaded regions denote areas where both inequalities hold true: $\{\mathbf{x} : \langle \mathbf{a}_{i_1}, \mathbf{x} \rangle \leq b_{i_1} \text{ and } \langle \mathbf{a}_{i_2}, \mathbf{x} \rangle \leq b_{i_2}\}$. Notice that in (b), when the angle between $\mathbf{x}_{j-1} - \mathbf{x}_j$ and $\mathbf{s} - \mathbf{x}_j$ is obtuse, the orthogonal projection of estimation \mathbf{x}_{j-1} onto their intersection is guaranteed to be closer to the solution set S . On the other hand, when that angle is acute we see exactly the opposite, as in (c). We can quantify this notion by the following definition.

Definition 3.1. For an $r \times d$ matrix \mathbf{A} and $\mathbf{b} \in \mathbb{R}^r$, for row i denote by \tilde{H}_i and H_i the half-space $\tilde{H}_i = \{\langle \mathbf{a}_i, \mathbf{x} \rangle \leq b_i\}$ and hyperplane $H_i = \{\langle \mathbf{a}_i, \mathbf{x} \rangle = b_i\}$, respectively, and write P_S as the orthogonal projection onto a convex set S . An *obtuse* (m, β) row paving of the matrix \mathbf{A} is an (m, β) row paving $T = \{\tau_1, \dots, \tau_m\}$ that also satisfies the following. Let $\tau \in T$ and let $\mathbf{s} \in \cap_{i \in \tau} \tilde{H}_i$, $\mathbf{w} \in \cap_{i \in \tau} \tilde{H}_i^c$, and $\mathbf{z} = P_{\cap_{i \in \tau} H_i} \mathbf{w}$. Then

$$\langle \mathbf{w} - \mathbf{z}, \mathbf{s} \rangle < 0.$$

In other words, the angle between $\mathbf{w} - \mathbf{z}$ and \mathbf{s} (and thus $\mathbf{s} - \mathbf{z}$) is obtuse.

We will see that performing block projections on the inequalities in the system only makes sense when one can obtain an obtuse row paving. We will use $\mathbf{w} = \mathbf{x}_{j-1}$, $\mathbf{z} = \mathbf{x}_j$, and $\mathbf{s} \in S$. Notice that if $i_1, i_2 \in \tau \in T$, then the partition used in the system depicted in Figure 1 (c) does not constitute an obtuse row paving.

We conduct two simple experiments to demonstrate the different behavior of the algorithm. In all cases the matrix \mathbf{A} is a 300×100 matrix with standard normal entries, 100 rows correspond to inequalities,

Algorithm 3.2 Double Block Kaczmarz Method for a System of Inequalities

Input:

- Matrix \mathbf{A} with dimension $n \times d$
- Right-hand side \mathbf{b} with dimension n
- Partition $T' = \{\tau'_1, \dots, \tau'_{m'}\}$ of the row indices $\{1, \dots, n_i\}$
- Partition $T = \{\tau_1, \dots, \tau_m\}$ of the row indices $\{1, \dots, n_e\}$
- Initial iterate \mathbf{x}_0 with dimension d
- Convergence tolerance $\varepsilon > 0$

Output: An estimate $\hat{\mathbf{x}}$ to the solution of the system (1.3)

```

j ← 0
repeat
  j ← j + 1
  Draw uniformly at random q from [0, 1]
  if q ≤  $\frac{\beta m}{\beta' m' + \beta m}$ 
    Choose a block  $\tau$  uniformly at random from T
     $\mathbf{x}_j \leftarrow \mathbf{x}_{j-1} + (\mathbf{A}_\tau)^\dagger (\mathbf{b}_\tau - \mathbf{A}_\tau \mathbf{x}_{j-1})$  (Solve least-squares approximation)
  else
    Choose a block  $\tau'$  uniformly at random from T'
    Set  $\sigma = \{i \in \tau' : \langle \mathbf{a}_i, \mathbf{x}_{j-1} \rangle > b_i\} \subset \tau'$  (Select unsatisfied subset)
     $\mathbf{x}_j \leftarrow \mathbf{x}_{j-1} + (\mathbf{A}_\sigma)^\dagger (\mathbf{b}_\sigma - \mathbf{A}_\sigma \mathbf{x}_{j-1})$  (Solve least-squares approximation)
until  $\|e(\mathbf{A}\mathbf{x}_j - \mathbf{b})\|_2^2 \leq \varepsilon^2$ 
 $\hat{\mathbf{x}} \leftarrow \mathbf{x}_j$ 

```

and \mathbf{b} is generated so that the solution set S is non-empty. We measure the residual error which we define as $\|e(\mathbf{A}\mathbf{x}_j - \mathbf{b})\|_2$. Figure 2 (a) shows the behavior of the block method with this matrix and a row paving obtained via a random row partition of 30 blocks (10 rows per block). This generation will create a matrix with paving that with very high probability is not an obtuse row paving. As Figure 2 demonstrates, the block method does not converge to a solution in this case. However, as Figure 2 (c) shows, the simple Kaczmarz method succeeds in identifying a point in the solution space. Next, we create a matrix in the exact same way, and create the same random row paving. Then, however, we iterate through every block in the paving corresponding to inequalities and if two rows i and k in a block satisfy $\langle \mathbf{a}_i, \mathbf{a}_k \rangle > 0$, we replace row \mathbf{a}_i with $-\mathbf{a}_i$ and entry b_i with $-b_i$. This guarantees every block in the paving yields a geometry like that shown in Figure 1 (b), and gives an obtuse row paving. Note that of course this changes the solution space as well so one cannot employ this strategy in general. We then add positive values to the entries in \mathbf{b} corresponding to inequalities to ensure the solution set S is non-empty. With this new system and paving, we again run the block method and see that the method now converges to a point in the solution set, as seen in Figure 2 (b).

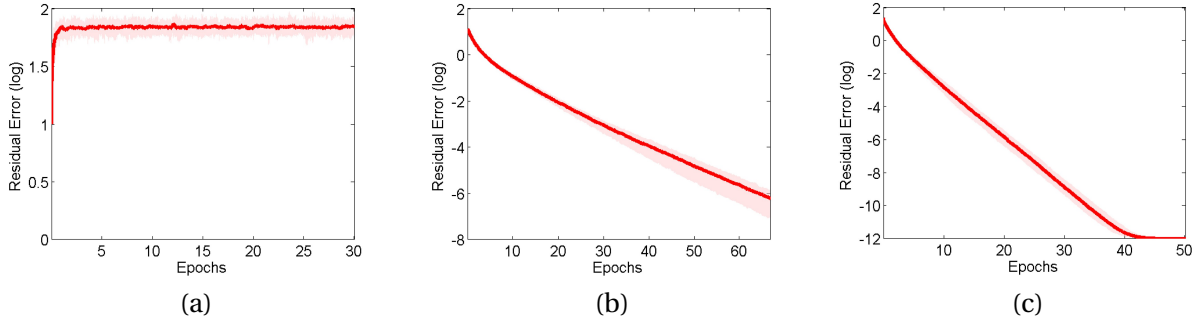


Figure 2 Residual error of the Kaczmarz Method per epoch: a) Median residual error of block method over 40 trials for matrix \mathbf{A} not having an obtuse row paving, b) Median residual error of block method over 40 trials for matrix \mathbf{A} using an obtuse row paving, c) Median residual error of simple method over 40 trials for same matrix as in a). Shaded region spans across minimum and maximum values over all trials and solid line denotes median value.

With this definition we obtain the following result, whose proof can be found in the appendix.

Theorem 3.2. *Let \mathbf{A} satisfy the assumptions of Theorem 2.1 and in addition have an obtuse (m', β') row paving of \mathbf{A}_\leq . Let x_1, \dots denote the iterates of Algorithm 3.2. Then using the notation of Theorem 2.1,*

$$\mathbb{E}[d(\mathbf{x}_j, S)^2] \leq \left[1 - \frac{1}{L^2(\beta' m' + \beta m)} \right]^j d(\mathbf{x}_0, S)^2.$$

Note that row pavings of standardized matrices can be obtained readily, often by random partitions [35, 36, 24], whereas obtuse row pavings may be much more challenging to obtain in general. Of course, by default the trivial paving which assigns each set τ to a single row always admits an obtuse row paving. We focus on Algorithm 2.1 which paves only \mathbf{A}_\leq , and leave further analysis of Algorithm 3.2 and constructions of obtuse row pavings for future work.

4. EXPERIMENTS

We use MATLAB to run some experiments using random matrices to test the convergence of the block Kaczmarz method applied to a system of equalities and inequalities. In each experiment, we create a random 500 by 100 matrix \mathbf{A} where each element is an independent standard normal random variable.

Each entry is then divided by the norm of its row so that the matrix is standardized. The first 400 rows of matrix \mathbf{A} compose $\mathbf{A}_=$, and the remaining 100 rows are set as inequalities of \mathbf{A}_\leq in the method described by (1.3). The experiments are run using the following procedure. For each of 100 trials,

- (1) Create matrix \mathbf{A} in the manner described above.
- (2) Create \mathbf{x}_\star where each entry is selected independently from a standard normal distribution. Set $\mathbf{b} = \mathbf{A}\mathbf{x}_\star$.
- (3) Pave submatrix $\mathbf{A}_=$ into 16 blocks with 25 equalities per block by a random partitioning of the rows.
- (4) Set initial approximations $\mathbf{x}_0^{\text{block}} = \mathbf{x}_0^{\text{simp}} = \mathbf{A}^* \mathbf{b}$.
- (5) Draw q uniformly at random from $[0, 1]$.
 - (a) If $q \leq \frac{n_e}{n}$, choose block $\{1, \dots, m\}$ uniformly at random and update iterate $\mathbf{x}_j^{\text{block}}$ using (1.9). (Note that the threshold $\frac{n_e}{n}$ is different than that given in the main algorithm and theorem, but it is easier to calculate and seems to work fine in practice.)
 - (b) Else, choose a row uniformly at random from $\{401, \dots, 500\}$ and update iterate $\mathbf{x}_j^{\text{block}}$ using (1.6).
 - (c) Update iterate $\mathbf{x}_j^{\text{simp}}$ using (1.6).

For both the simple and block algorithms, the median, minimum, and maximum values of the residual $\|e(\mathbf{A}\mathbf{x}_j - \mathbf{b})\|_2^2$ of the 100 trials are recorded for each iteration j .

Figure 3 compares the performance of the block Kaczmarz method used in this paper and the standard Kaczmarz method described by Leventhal and Lewis [19]. The plot in Figure 3 (a) compares convergence per iteration. As the block Kaczmarz method enforces multiple equalities per iteration, it is unsurprising that it performs better in this experiment. Figure 3 (b) displays the convergence of the two methods per epoch. The block Kaczmarz algorithm has an epoch of $m + n_i$ iterations, and the standard Kaczmarz method has an epoch of size n . Here, to be fair we only count an iteration towards an epoch if the estimated solution $\mathbf{x}_j \neq \mathbf{x}_{j-1}$. Thus in the case where a chosen inequality is already satisfied for iteration j , this iteration does not count towards an epoch since no computation is being performed. We noticed, however, that whether or not we modified the count in this way, the behavior still produces results very similar to Figure 3. Once again the experiments yielded faster convergence with the block Kaczmarz approach. It is interesting to compare the results of Figure 3 (b) and those of Figure 2 (b) and (c). The per-epoch convergence of the methods and whether the block or standard appears faster varies slightly and depends on both the number of rows and columns. In general, the per-epoch convergence rates are reasonably comparable, as the analysis suggests. However, Figure 3 (c) compares the rate of convergence of the two algorithms by plotting the residual against the CPU time expended in the simulation. We believe that the ability to utilize efficient matrix–vector multiplication gives the method significantly improved convergence per second relative to the standard Kaczmarz algorithm, although other mechanisms may certainly be at work as well.

5. CONCLUSION AND RELATED WORK

The Kaczmarz algorithm was first proposed in [18]. Kaczmarz demonstrated that the method converged to the solution of linear system $\mathbf{Ax} = \mathbf{b}$ for square, non-singular matrix \mathbf{A} . Since then, the method has been utilized in the context of computer tomography as the *Algebraic Reconstruction Technique* (ART) [12, 3, 21, 16]. Empirical results suggested that randomized selection offered improved convergence over the cyclic scheme [13, 15]. Strohmer and Vershynin [31] were the first to prove an expected linear convergence rate using a randomized Kaczmarz algorithm with specific random control. This result was extended by Needell [22] to apply to inconsistent systems, which shows a linear convergence rate to within a fixed radius around the least-squares solution. Almost-sure convergence guarantees were recently proved by Chen and Powell [6]. Zouzias and Freris [41] analyze a modified version of the method

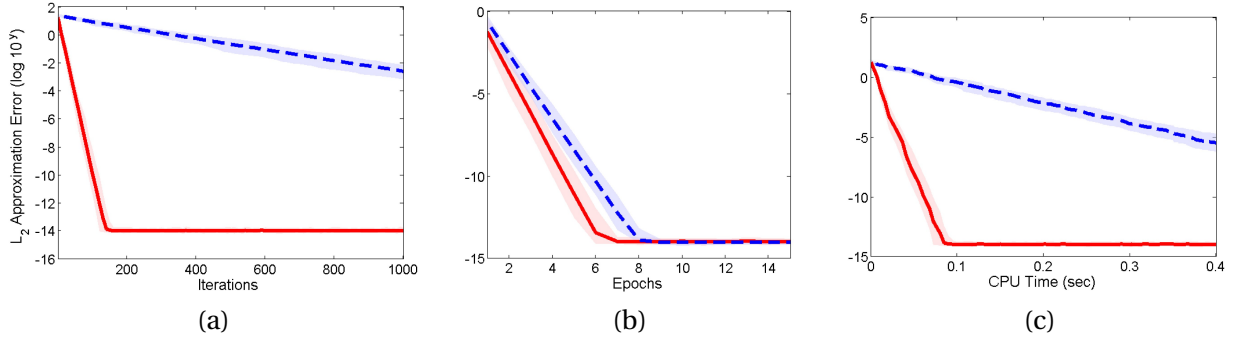


Figure 3 Residual error of the Block Kaczmarz Method (solid red) vs. Simple Kaczmarz Method (dashed blue) as a function of (a) Iterations, (b) Epochs, (c) CPU time. Shaded region spans from minimum to maximum value over 100 trials; lines denote the median value.

in the inconsistent case, using a variant motivated by Popa [26] to reduce the residual and thereby converge to the least squares solution. Relaxation parameters can also be introduced to obtain convergence to the least squares solution, see e.g. [39, 4, 32, 14], and partially weighted sampling can lead to a tradeoff between convergence rate and radius [23]. Liu, Wright, and Sridhar [20] discuss applying a parallelized variant of the randomized Kaczmarz method, demonstrating that the convergence rate can be increased almost linearly by bounding the number of processors by a multiple of the number of rows of A .

The block Kaczmarz updating method was introduced by Elfving [9] as a special case of the more general framework by Eggermont et.al. [8]. The notion of using blocking in projection methods is certainly not new, and there is a large amount of literature on these types of methods, see e.g. [40, 3] and references therein. Needell and Tropp [24] provide the first analysis showing an expected linear convergence rate which depends on the properties of the matrix A and of the submatrices A_{τ} resulting from the paving, connecting pavings and the block Kaczmarz scheme. The use of specialized blocks appears elsewhere, in particular, the works of Popa use blocks with orthogonal rows that are beneficial for the block Kaczmarz method [26, 27, 28]. Needell, Zhao, and Zouzias [25] expand on the results from [24] and [41] to demonstrate convergence to the least-squares solution for an inconsistent system using the block Kaczmarz method. Again the block approach can yield faster convergence than the simple method.

The Kaczmarz method was first applied to a system of equalities and inequalities by Leventhal and Lewis [19], who also consider polynomial constraints with the method. They give a linear convergence rate to the feasible solution space S , using $\|A\|_F^2$ and the Hoffman constant [17]. We apply the block Kaczmarz scheme to the system described in [19], combining their result with that of Needell and Tropp [24] to acquire a completely generalized result. We highlight several important complications which arise when attempting to apply the block scheme to inequalities. Nonetheless, whether a paving is used only partially or for the complete system, significant reduction in computational time can be achieved.

5.1. Future Work. There are many interesting open problems related to the block Kaczmarz method and linear systems with inequalities. It has been well observed in the literature that selecting rows (or blocks) *without replacement* rather than with replacement as in the theoretical results leads to faster a convergence rate empirically [29, 24]. When selecting without replacement, independence between iterations vanishes, making a theoretical analysis more challenging. Secondly, it would be interesting to further investigate the use of obtuse row pavings. In systems with a large number of inequalities, the ability to pave the submatrix A_{\leq} with an obtuse row paving would lead to significantly faster convergence. In that case, one may like to identify a more general geometric property about the system that permits such pavings or an alternative formulation that offers convergence of the full block method.

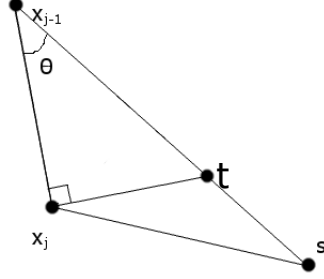


Figure 4 Geometry of system.

APPENDIX A. PROOF OF THEOREM 3.2

Proof. Fix an iteration j of Algorithm 3.2. As in the proof of Theorem 2.1, if a block of equalities is selected this iteration, then we again have (2.1). So we next instead consider the case when a block of inequalities is selected, and call this block τ' , and its pruned subset σ . Set $\mathbf{s} = P_S \mathbf{x}_{j-1}$, where again P_S denotes the orthogonal projection onto the solution set S . If we write $\tilde{H}_i = \{\mathbf{x} : \langle \mathbf{a}_i, \mathbf{x} \rangle \leq b_i\}$ and $H_i = \{\mathbf{x} : \langle \mathbf{a}_i, \mathbf{x} \rangle = b_i\}$, then by their definitions we have

$$\mathbf{s} \in \cap_{i \in \sigma} \tilde{H}_i, \quad \mathbf{x}_{j-1} \in \cap_{i \in \sigma} \tilde{H}_i^c, \quad \text{and} \quad \mathbf{x}_j = P_{\cap_{i \in \sigma} H_i} \mathbf{x}_{j-1}.$$

Then since σ is part of an obtuse paving, the angle between $\mathbf{x}_j - \mathbf{x}_{j-1}$ and $\mathbf{s} - \mathbf{x}_{j-1}$ must be obtuse. There thus exists a point \mathbf{t} on the line segment $L = \{\gamma \mathbf{x}_{j-1} + (1 - \gamma) \mathbf{s} : 0 \leq \gamma \leq 1\}$ such that $\mathbf{x}_{j-1} - \mathbf{x}_j$ and $\mathbf{t} - \mathbf{x}_j$ are orthogonal (see Figure 4).

Now since $\mathbf{t} \in L$, we have $\|\mathbf{t} - \mathbf{x}_{j-1}\|_2 \leq \|\mathbf{x}_{j-1} - \mathbf{s}\|_2$, and thus letting θ denote the angle between $\mathbf{x}_j - \mathbf{x}_{j-1}$ and $\mathbf{t} - \mathbf{x}_{j-1}$ (see Figure 4), we have

$$\begin{aligned} \|\mathbf{x}_j - \mathbf{x}_{j-1}\|_2 &\leq \|\mathbf{x}_{j-1} - \mathbf{s}\|_2 \cdot \frac{\|\mathbf{x}_j - \mathbf{x}_{j-1}\|_2}{\|\mathbf{t} - \mathbf{x}_{j-1}\|_2} \\ &= \|\mathbf{x}_{j-1} - \mathbf{s}\|_2 \cdot \cos \theta \\ &= \frac{\|\mathbf{x}_j - \mathbf{x}_{j-1}\|_2 \cdot \|\mathbf{s} - \mathbf{x}_{j-1}\|_2 \cdot \cos \theta}{\|\mathbf{x}_j - \mathbf{x}_{j-1}\|_2} \\ &= \frac{\langle \mathbf{s} - \mathbf{x}_{j-1}, \mathbf{x}_j - \mathbf{x}_{j-1} \rangle}{\|\mathbf{x}_j - \mathbf{x}_{j-1}\|_2} \\ &= \frac{-\langle \mathbf{x}_{j-1} - \mathbf{s}, \mathbf{x}_j - \mathbf{x}_{j-1} \rangle}{\|\mathbf{x}_j - \mathbf{x}_{j-1}\|_2}. \end{aligned}$$

Thus, we have that

$$\langle \mathbf{x}_{j-1} - \mathbf{s}, \mathbf{x}_j - \mathbf{x}_{j-1} \rangle \leq -\|\mathbf{x}_j - \mathbf{x}_{j-1}\|_2^2.$$

By the definition of \mathbf{x}_j , this means that

$$\langle \mathbf{x}_{j-1} - \mathbf{s}, \mathbf{A}_\sigma^\dagger (\mathbf{b}_\sigma - \mathbf{A}_\sigma \mathbf{x}_{j-1}) \rangle \leq -\|\mathbf{A}_\sigma^\dagger (\mathbf{b}_\sigma - \mathbf{A}_\sigma \mathbf{x}_{j-1})\|_2^2. \quad (\text{A.1})$$

Using this along with the paving properties we see that

$$\begin{aligned}
\|\mathbf{x}_j - \mathbf{s}\|_2^2 &= \|\mathbf{x}_{j-1} - \mathbf{s} + \mathbf{A}_\sigma^\dagger(\mathbf{b}_\sigma - \mathbf{A}_\sigma \mathbf{x}_{j-1})\|_2^2 \\
&= \|\mathbf{x}_{j-1} - \mathbf{s}\|_2^2 + 2\langle \mathbf{x}_{j-1} - \mathbf{s}, \mathbf{A}_\sigma^\dagger(\mathbf{b}_\sigma - \mathbf{A}_\sigma \mathbf{x}_{j-1}) \rangle \\
&\quad + \|\mathbf{A}_\sigma^\dagger(\mathbf{b}_\sigma - \mathbf{A}_\sigma \mathbf{x}_{j-1})\|_2^2 \\
&\leq \|\mathbf{x}_{j-1} - \mathbf{s}\|_2^2 - \|\mathbf{A}_\sigma^\dagger(\mathbf{b}_\sigma - \mathbf{A}_\sigma \mathbf{x}_{j-1})\|_2^2 \\
&= d(\mathbf{x}_{j-1}, S)^2 - \|\mathbf{A}_\sigma^\dagger(\mathbf{b}_\sigma - \mathbf{A}_\sigma \mathbf{x}_{j-1})\|_2^2 \\
&\leq d(\mathbf{x}_{j-1}, S)^2 - \frac{1}{\beta'} \|\mathbf{b}_\sigma - \mathbf{A}_\sigma \mathbf{x}_{j-1}\|_2^2.
\end{aligned}$$

Thus, taking expectation (over the choice of τ' , conditioned on previous choices), yields

$$\begin{aligned}
\mathbb{E}[d(\mathbf{x}_j, S)^2] &\leq \mathbb{E}\|\mathbf{x}_j - \mathbf{s}\|_2^2 \\
&\leq d(\mathbf{x}_{j-1}, S)^2 - \frac{1}{\beta'} \mathbb{E}\|\mathbf{b}_\sigma - \mathbf{A}_\sigma \mathbf{x}_{j-1}\|_2^2 \\
&= d(\mathbf{x}_{j-1}, S)^2 - \frac{1}{\beta'} \mathbb{E}\|e(\mathbf{b}_{\tau'} - \mathbf{A}_{\tau'} \mathbf{x}_{j-1})\|_2^2 \\
&= d(\mathbf{x}_{j-1}, S)^2 - \frac{1}{m' \beta'} \sum_{\tau' \in T'} \|e(\mathbf{b}_{\tau'} - \mathbf{A}_{\tau'} \mathbf{x}_{j-1})\|_2^2 \\
&= d(\mathbf{x}_{j-1}, S)^2 - \frac{1}{m' \beta'} \|e(\mathbf{b}_\leq - \mathbf{A}_\leq \mathbf{x}_{j-1})\|_2^2.
\end{aligned}$$

Combining this with (2.1) and letting E_+ and E_\leq denote the events that a block from T and a block from T' is selected, respectively, we have

$$\begin{aligned}
\mathbb{E}[(d(\mathbf{x}_j, S)^2)] &= p \cdot \mathbb{E}[d(\mathbf{x}_j, S)^2 | E_+] + (1-p) \cdot \mathbb{E}[d(\mathbf{x}_j, S)^2 | E_\leq] \\
&\leq p \left[d(\mathbf{x}_{j-1}, S)^2 - \frac{1}{\beta m} \sum_{i \in I_+} e(\mathbf{A}_+ \mathbf{x}_{j-1} - \mathbf{b}_+)_i^2 \right] \\
&\quad + (1-p) \left[d(\mathbf{x}_{j-1}, S)^2 - \frac{1}{m' \beta'} \|e(\mathbf{b}_\leq - \mathbf{A}_\leq \mathbf{x}_{j-1})\|_2^2 \right] \\
&= d(\mathbf{x}_{j-1}, S)^2 - p \cdot \frac{1}{\beta m} \sum_{i \in I_+} e(\mathbf{A}_+ \mathbf{x}_{j-1} - \mathbf{b}_+)_i^2 \\
&\quad - (1-p) \cdot \frac{1}{m' \beta'} \|e(\mathbf{b}_\leq - \mathbf{A}_\leq \mathbf{x}_{j-1})\|_2^2
\end{aligned}$$

Since $p = \frac{\beta m}{\beta' m' + \beta m}$, we have $\frac{1-p}{\beta' m'} = \frac{1}{\beta' m' + \beta m}$ and we can simplify

$$\begin{aligned}
\mathbb{E} [d(\mathbf{x}_j, S)^2] &\leq d(\mathbf{x}_{j-1}, S)^2 - \frac{1}{\beta' m' + \beta m} \left[\sum_{i \in I_-} e(\mathbf{A}_- \mathbf{x}_{j-1} - \mathbf{b}_-)_i^2 \right. \\
&\quad \left. + \|e(\mathbf{b}_- - \mathbf{A}_- \mathbf{x}_{j-1})\|_2^2 \right] \\
&= d(\mathbf{x}_{j-1}, S)^2 - \frac{1}{\beta' m' + \beta m} \|e(\mathbf{A} \mathbf{x}_{j-1} - \mathbf{b})\|_2^2 \\
&\leq d(\mathbf{x}_{j-1}, S)^2 - \frac{1}{L^2(\beta' m' + \beta m)} \cdot d(\mathbf{x}_{j-1}, S)^2 \\
&= \left[1 - \frac{1}{L^2(\beta' m' + \beta m)} \right] d(\mathbf{x}_{j-1}, S)^2,
\end{aligned}$$

where we have utilized the Hoffman bound (1.7) in the second inequality.

Iterating this relation along with independence of the random control completes the proof. \square

REFERENCES

- [1] Bourgain, J., Tzafriri, L.: Invertibility of “large” submatrices with applications to the geometry of Banach spaces and harmonic analysis. *Israel J. Math.* **57**(2), 137–224 (1987). DOI 10.1007/BF02772174. URL <http://dx.doi.org/10.1007/BF02772174>
- [2] Bourgain, J., Tzafriri, L.: On a problem of Kadison and Singer. *J. Reine Angew. Math.* **420**, 1–43 (1991). JRMMA8; 46L05 (46L30 47B35 47D25); 1124564 (92j:46104); H. Halpern
- [3] Byrne, C.L.: *Applied iterative methods*. A K Peters Ltd., Wellesley, MA (2008)
- [4] Censor, Y., Eggermont, P.P.B., Gordon, D.: Strong underrelaxation in kaczmarz’s method for inconsistent systems. *Numer. Math.* **41**(1), 83–92 (1983)
- [5] Censor, Y.: Row-action methods for huge and sparse systems and their applications. *SIAM Review* **23**(4), 444–466 (1981)
- [6] Chen, X., Powell, A.: Almost sure convergence of the Kaczmarz algorithm with random measurements. *J. Fourier Anal. Appl.* pp. 1–20 (2012). URL <http://dx.doi.org/10.1007/s00041-012-9237-2>. DOI 10.1007/s00041-012-9237-2
- [7] Chrétien, S., Darses, S.: Invertibility of random submatrices via tail decoupling and a matrix chernoff inequality. *Statist. Probab. Lett.* **82**(7), 1479–1487 (2012)
- [8] Eggermont, P.P.B., Herman, G.T., Lent, A.: Iterative algorithms for large partitioned linear systems, with applications to image reconstruction. *Linear Algebra Appl.* **40**, 37–67 (1981). DOI 10.1016/0024-3795(81)90139-7. URL [http://dx.doi.org/10.1016/0024-3795\(81\)90139-7](http://dx.doi.org/10.1016/0024-3795(81)90139-7)
- [9] Elfving, T.: Block-iterative methods for consistent and inconsistent linear equations. *Numer. Math.* **35**(1), 1–12 (1980). DOI 10.1007/BF01396365. URL <http://dx.doi.org/10.1007/BF01396365>
- [10] Feichtinger, H.G., Cenkler, C., Mayer, M., Steier, H., Strohmer, T.: New variants of the POCS method using affine subspaces of finite codimension with applications to irregular sampling. In: *Applications in Optical Science and Engineering*, pp. 299–310. International Society for Optics and Photonics (1992)
- [11] Feichtinger, H.G., Strohmer, T.: A kaczmarz-based approach to nonperiodic sampling on unions of rectangular lattices. In: *SampTA ’95: 1995 Workshop on Sampling Theory and Applications*, pp. 32–37. Jurmala, Latvia (1995)
- [12] Gordon, R., Bender, R., Herman, G.T.: Algebraic reconstruction techniques (ART) for three-dimensional electron microscopy and X-ray photography. *J. Theoret. Biol.* **29**, 471–481 (1970)
- [13] Hamaker, C., Solmon, D.C.: The angles between the null spaces of X-rays. *J. Math. Anal. Appl.* **62**(1), 1–23 (1978)
- [14] Hanke, M., Niethammer, W.: On the acceleration of kaczmarz’s method for inconsistent linear systems. *Linear Algebra Appl.* **130**, 83–98 (1990)
- [15] Herman, G., Meyer, L.: Algebraic reconstruction techniques can be made computationally efficient. *IEEE T. Med. Imaging* **12**(3), 600–609 (1993)
- [16] Herman, G.T.: *Fundamentals of computerized tomography: image reconstruction from projections*. Springer (2009)
- [17] Hoffman, A.J.: On approximate solutions of systems of linear inequalities. *J. Research Nat. Bur. Standards* **49**, 263–265 (1952)
- [18] Kaczmarz, S.: Angenäherte auflösung von systemen linearer gleichungen. *Bull. Int. Acad. Polon. Sci. Lett. Ser. A* pp. 335–357 (1937)
- [19] Leventhal, D., Lewis, A.S.: Randomized methods for linear constraints: convergence rates and conditioning. *Math. Oper. Res.* **35**(3), 641–654 (2010). DOI 10.1287/moor.1100.0456. URL <http://dx.doi.org/10.1287/moor.1100.0456>

- [20] Liu, J., Wright, S.J., Srikrishna, S.: An asynchronous parallel randomized kaczmarz algorithm (2014). Available at [arXiv: 1401.4780](https://arxiv.org/abs/1401.4780)
- [21] Natterer, E.: The mathematics of computerized tomography, *Classics in Applied Mathematics*, vol. 32. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA (2001). DOI 10.1137/1.9780898719284. URL <http://dx.doi.org/10.1137/1.9780898719284>. Reprint of the 1986 original
- [22] Needell, D.: Randomized Kaczmarz solver for noisy linear systems. *BIT* **50**(2), 395–403 (2010). DOI 10.1007/s10543-010-0265-5. URL <http://dx.doi.org/10.1007/s10543-010-0265-5>
- [23] Needell, D., Srebro, N., Ward, R.: Stochastic gradient descent and the randomized kaczmarz algorithm (2013). Submitted
- [24] Needell, D., Tropp, J.A.: Paved with good intentions: Analysis of a randomized block Kaczmarz method. *Linear Algebra Appl.* **441**, 199–221 (2014)
- [25] Needell, D., Zhao, R., Zouzias, A.: Randomized block kaczmarz method with projection for solving least squares (2014). Submitted
- [26] Popa, C.: Block-projections algorithms with blocks containing mutually orthogonal rows and columns. *BIT* **39**(2), 323–338 (1999). DOI 10.1023/A:1022398014630. URL <http://dx.doi.org/10.1023/A:1022398014630>
- [27] Popa, C.: A fast Kaczmarz-Kovarik algorithm for consistent least-squares problems. *Korean J. Comput. Appl. Math.* **8**(1), 9–26 (2001)
- [28] Popa, C.: A Kaczmarz-Kovarik algorithm for symmetric ill-conditioned matrices. *An. Ştiinţ. Univ. Ovidius Constanţa Ser. Mat.* **12**(2), 135–146 (2004)
- [29] Recht, B., Ré, C.: Beneath the valley of the noncommutative arithmetic–geometric mean inequality: Conjectures, case studies, and consequences. In: *Proc. 25th Ann. Conf. Learning Theory*. Edinburgh (2012)
- [30] Sezan, M.I., Stark, H.: Applications of convex projection theory to image recovery in tomography and related areas. *Image Recovery: Theory and Application* pp. 155–270 (1987)
- [31] Strohmer, T., Vershynin, R.: A randomized Kaczmarz algorithm with exponential convergence. *J. Fourier Anal. Appl.* **15**(2), 262–278 (2009). DOI 10.1007/s00041-008-9030-4. URL <http://dx.doi.org/10.1007/s00041-008-9030-4>
- [32] Tanabe, K.: Projection method for solving a singular system of linear equations and its applications. *Numer. Math.* **17**(3), 203–214 (1971)
- [33] Tropp, J.A.: The random paving property for uniformly bounded matrices. *Studia Math.* **185**(1), 67–82 (2008). URL <http://dx.doi.org/10.4064/sm185-1-4.46B09> (15A52 46B20 60E15); 2379999 (2008k:46030); Sasha Sodin
- [34] Tropp, J.A.: Column subset selection, matrix factorization, and eigenvalue optimization. In: *Proc. Twentieth Annual ACM-SIAM Symposium on Discrete Algorithms*, pp. 978–986. SIAM, Philadelphia, PA (2009)
- [35] Tropp, J.A.: Improved analysis of the subsampled randomized hadamard transform. *Advances in Adaptive Data Analysis* **3**(01n02), 115–126 (2011)
- [36] Tropp, J.A.: User-friendly tail bounds for sums of random matrices. *Found. Comput. Math.* **12**(4), 389–434 (2012)
- [37] Vershynin, R.: John’s decompositions: Selecting a large part, *Israel Journal of Mathematics*, **122**(1), 253–277 (2001). *Inst. Math. Statist*, Beachwood, OH (2006). URL <http://dx.doi.org/10.1214/074921706000000815>. 46B09 (46B07 46B20); 2387766 (2009h:46023); Dirk Werner
- [38] Vershynin, R.: Random sets of isomorphism of linear operators on Hilbert space, *High dimensional probability*, vol. 51, pp. 148–154. *Inst. Math. Statist*, Beachwood, OH (2006). URL <http://dx.doi.org/10.1214/074921706000000815>. 46B09 (46B07 46B20); 2387766 (2009h:46023); Dirk Werner
- [39] Whitney, T.M., Meany, R.K.: Two algorithms related to the method of steepest descent. *SIAM J. Numer. Anal.* **4**(1), 109–118 (1967)
- [40] Xu, J., Zikatanov, L.: The method of alternating projections and the method of subspace corrections in Hilbert space. *J. Amer. Math. Soc.* **15**(3), 573–597 (2002). DOI 10.1090/S0894-0347-02-00398-3. URL <http://dx.doi.org/10.1090/S0894-0347-02-00398-3>
- [41] Zouzias, A., Freris, N.M.: Randomized extended kaczmarz for solving least squares. *SIAM J. Matrix Anal. A.* **34**(2), 773–793 (2013)